

# Reading Comprehension using Entity-based Memory Networks

Xun Wang<sup>1,2</sup>, Katsuhito Sudoh<sup>1</sup>, Masaaki Nagata<sup>1</sup>,  
Tomohide Shibata<sup>2</sup>, Daisuke Kawahara<sup>2</sup>, and Sadao Kurohashi<sup>2</sup>

<sup>1</sup> NTT Communication Science Laboratories, Kyoto, Japan,  
wang.xun, sudoh.katsuhito, nagata.masaaki@lab.ntt.co.jp

<sup>2</sup> Kyoto University, Kyoto, Japan  
shibata,dk,kuro@i.kyoto-u.ac.jp

**Abstract.** This paper introduces a novel neural network model for question answering, the *entity-based memory networks*. It enhances neural networks' ability of representing and calculating information over a long period by keeping records of entities contained in text. The core component is a memory pool which comprises entities' states. These entities' states are continuously updated according to the input text. Questions with regard to the input text are used to search the memory pool for related entities and answers are further predicted based on the states of retrieved entities. Compared with previous memory network models, the proposed model is capable of handling fine-grained information and more sophisticated relations based on entities. We formulated several different tasks as question answering problems and tested the proposed model. Experiments reported satisfying results.

**Keywords:** Text Comprehension, Entity Memory Network, Question Answering

## 1 Introduction

It has long been a major concern of the natural language processing (NLP) community to enable computers to understand text as humans do. A lot of NLP tasks have been tensely studied towards this goal such as information retrieval, semantical role labelling, textual entailment and so on. Among them, questions answering is of great importance and has been a huge challenge. A question answering (QA) task is to predict an answer for a given question with regard to related information. It can be formulated as a map  $f : \{related\_text, question\} \rightarrow \{answer\}$  [12]. To predict the correct answer, computers are firstly required to “understand” the text.

Shallow features such as bag-of-words, token frequencies and so on are unable to capture the rich information in text. Often outside knowledge is required towards better performances. Traditional approaches heavily rely on rules or structured knowledge developed by experts or crowd sourcing [21,18]. Relational databases constructed from predicate argument triples also serve as a source of knowledge [14,23]. Problems with these approaches lie in at least two aspects. Firstly the construction of structured knowledge is both time and money consuming. Secondly it is a huge challenge to design

models flexible and powerful enough to learn to employ these information extracted [6]. Thus the progress of using machine learning for QA has been slow.

Recently the emergence of deep neural networks and distributed representations sheds light on such methods. Representing all the features using vectors provides a unified representational form for all the necessary information. Outside knowledge learnt from large corpus can be encoded into word vectors. Information obtained locally is also represented using vectors. Deep neural network models with many layers are designed to fuse information obtained from different sources [4,7].

A notable breakthrough is to employ memories in neural networks. The representative model is named the memory network [28]. The key of memory network is to store historical sentences in a memory pool. The model is trained to look for related sentences when a question comes. Then based on the related sentences, an answer is predicted for the question. Memory network remembers all sentences it has read so that it can look for useful ones when facing questions. This model and its variants have been proved useful in a series of tasks [28,25,1].

One problem with memory networks is that using sentence vectors as elementary units of information makes it difficult to fully explore the information contained in text. Often is the case that in a long sentence, only part of the sentence is related to the questions. Therefore taking the whole sentence into consideration makes it hard to focus on the information that are related to questions. Besides, learning sentence representations itself is a growing field.

We propose to focus on entities rather than sentences. Entities refer to anything that exist in reality or are purely hypothetical. We assume that text can be projected to a world of entities. The key of conducting comprehension and reasoning over text is to identify its containing entities and analyze the states of these entities and the relations between them. We keep a memory pool of entities and use the input sentences to update the states of these entities. Questions are answered based on the states of related entities. The proposed model deals with fine-grained information by using entities. The introduced model is named as *entity-based memory network*. It is tested on several datasets, including the toy bAbI dataset [27], large movie review dataset [15] and the machine comprehension test dataset [20]. Results show we have achieved satisfying results using the entity-based memory network. The rest of the paper is organized as follows: Section 2 reviews some previous work. Section 3 describes our approaches and elaborates the details. Section 4 presents the experiments and the analysis. Section 5 concludes the paper.

## 2 Memories in Deep Neural Networks

QA has a long history and a lot of methods have been developed to address this problem [19,14,2]. Recently the development of neural models leads to series of work on question answering [5,4,7]. Among them closely related to our work is the Memory Network (MNN) [28]. The memory network contains four parts: the input module which converts sentences into vectors, the memory which keeps all sentence vectors a retrieval module and a response module. Whenever a question comes, the question is turned into a vector and the question vector is used to retrieve the memory for related sentences. The

response module is used to predict an answer based on the related sentences. The core component is the memory pool that stores all the input sentences so that they can be retrieved later to answer questions. This model contains several neural networks which are jointly optimized according to the task. Experiments on a toy dataset show that this model is able to answer simple questions according to the input text. Fig. 1(a) illustrates the memory network.

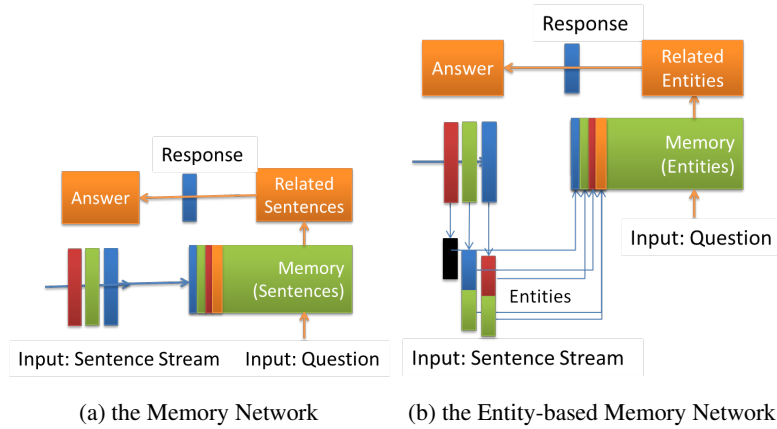


Fig. 1: Comparison of the memory network and the proposed entity-based memory network model. Sentences are decomposed into entities and then stored in the memory for later retrieval.

Later [10] propose the Dynamic Memory Network (DMNN) which introduces the attention mechanism into the memory network model. When retrieving memories, the location of the next related sentence is predicted according to the related sentences identified in the previous iterations. Using the attention mechanism, they obtain further improvements. Some other work [25,1] propose other variants of MNN by introducing additional memory network modules. These work focuses on storing sentence vectors for later retrieval with no exceptions. Most of them have been tested on the toy dataset bAbI [27] and are reported to have achieved satisfying results. When further tested on some practical tasks, these models also show the ability to produce results as good as existing state-of-the-art systems or even better results. Memory networks store sentence vectors as memories and have the superiority of processing information from a large scale. Experiment results they reported on a series of tasks are concrete proofs.

But there is also a problem with the memory networks as we have stated. Taking sentence vectors as input means that it is difficult to further analyze and take advantages of relations between smaller text units, such as entities. For example, when an entity  $e_a$  of sentence  $A$  interacts with another entity  $e_b$  of sentence  $B$ , we have to take the whole sentences  $A$  and  $B$  into consideration rather than just focus on  $e_a$  and  $e_b$ . This inevitably brings about noise and damages the comprehension of text. The failure of obtaining

fine-grained information prevents any further improvements. In the proposed entity-based model, we focus on entities directly and avoid bringing in redundant information.

### 3 Approaches

#### 3.1 Overview

Firstly we use an example to illustrate how the model works. Below we show a piece of text which contains 4 sentences and 2 questions. There are 7 entities in total, all of them underlined.

- 1) Mary moved to the bathroom. 2) John went to the hallway. 3) Where is Mary? Bathroom.  
4) Daniel went back to the hallway. 5) Sandra moved to the garden. 6) Where is Daniel? Hallway.

Fig. 2: An Example from bAbI, a toy dataset for question answering [27].

This text is elaborated around the 7 entities. It describes how their states change (i.e., the change of a character’s location) when the story goes on. Note that here all the entities are concrete concepts that exist in reality. It is also possible to talk about abstract concepts.

The core of the proposed model are entities. We take Sentence 1 ( $S_1$ ) as input and extract the entities it contains {Mary, bathroom}. Vectors representing the states of these entities are initialized using some pre-learned word embeddings  $\{\vec{Mary}, \vec{bathroom}\}$  and stored in a memory pool. Meanwhile, we turn  $S_1$  into a vector ( $\vec{S_1}$ ) using an autoencoder model<sup>3</sup>. Then we use the sentence vector  $\vec{S_1}$  to update the entities’ states  $\{\vec{Mary}, \vec{bathroom}\}$ . The goal is to reconstruct  $\vec{S_1}$  solely from  $\{\vec{Mary}, \vec{bathroom}\}$ . In the same way, we process the following text ( $S_2$ ) and its containing entities (John, hallway) until encounter a question ( $S_3$ ).  $S_3$  is converted into a vector ( $\vec{S_3}$ ) following the same method that processes previous input text. Then taking  $\vec{S_3}$  as input, we retrieve related entities from the memory which now stores all the entities (Mary, bathroom, John, hallway) that appear before  $S_3$ . The related entities’ states are then used to produce a feature vector. In this case, (Mary and bathroom) are related to the question and their states are used for constructing the feature vector. Note the current states of the two entities (Mary and bathroom) are different from their initial values due to  $S_1$ . Based on the feature vector, we then use another neural network model to predict the answer to  $S_3$ .

The model monitors the entities involved in text and keeps updating their states according to the input. Whenever we have a question with regard to the text, we check the states of entities and predict an answer accordingly. The proposed model comprises of 4 modules, as is shown in Fig. 3. Each module is designed for a unique purpose and together they construct the *entity-based memory network* model.

<sup>3</sup> Note that the sentence vector is not used to answer question directly and it is also plausible to use other models to learn sentence representation.

1. **I: Input module.** Take as input a sentence and turn it into a vector. Meanwhile, extract all the entities it contains. The question is also processed using this module.
2. **G: Generalization module.** Update the states of related entities according to the input. For entities that are not contained in the memory pool, create a new memory slot for each of them and initialize these slots using pre-learned word embeddings.
3. **O: Output feature module.** It is triggered whenever a question arrives. Retrieve related entities according to the input question and then produce an output feature vector accordingly.
4. **R: Response module.** Generate the response according to the output feature vector.

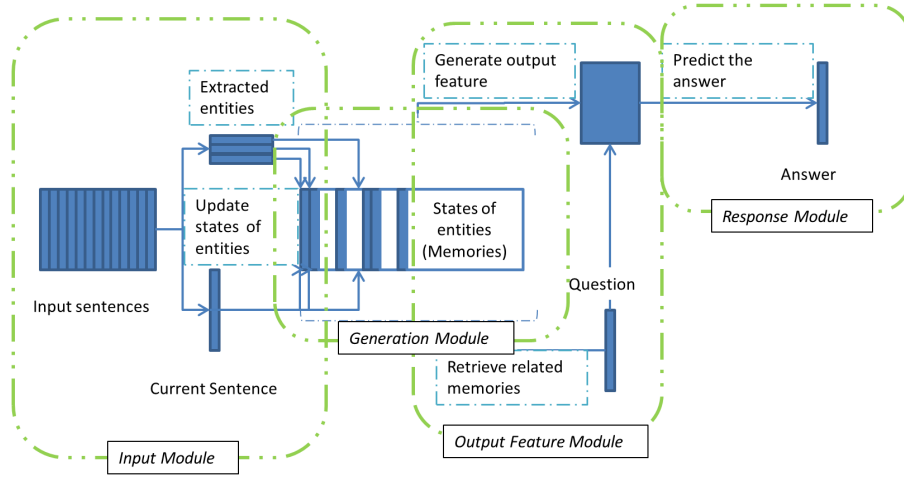


Fig. 3: Architecture of the entity-based memory network.  
The model is divided into four modules which are shown in the figure using squares.

### 3.2 Entity-based Memory Network Model

Here we present a formal description of the proposed model. Assume we have sentences  $S_1, S_2, \dots, S_n$  whose entities are annotated in advance as  $e_1, e_2, \dots, e_m$ .

*Input Module* We firstly turn each sentence  $S_i$  into its vector representation:

$$S_i = f_1(S_i) \quad (1)$$

*Generalization Module* For a sentence  $S_i$ , we collect all the entities it contains  $\{e_1^i, \dots, e_k^i, \dots, e_j^i\}$ . These entities' states  $\{e_k^i\}$  are simultaneously updated according to  $S_i$  as follows:

$$\{e_k^i\} = \arg \min_{\{e_k^i\}} (|S'_i - S_i|); S'_i = f_2(e_1^i, \dots, e_k^i, \dots, e_j^i); \quad (2)$$

$f_2$  is to reconstruct  $S_i$  only using the states of  $S_i$ 's containing entities  $\{e_k^i\}$ .  $\{e_k^i\}$  are updated to minimize the difference between  $S'_i$  and  $S_i$ . Recall that  $S_i$  is generated using  $f_1$  with the whole sentence  $S_i$  as input. We compress the information carried by  $S_i$  into a vector  $S_i$  and then unfold it into  $\{e_k^i\}$ .

After processing these sentences, we construct a memory pool which consists of entities whose states are regarded as capable of representing the information carried by the input text.

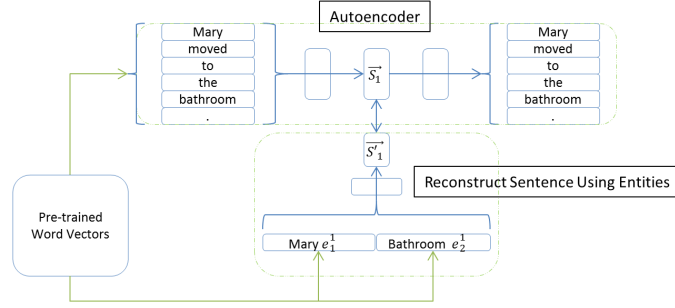


Fig. 4: the Generalization Module: Using  $S_1$  as an example, the autoencoder is used to convert the sentence into a vector  $S_1$  and the entities contained in  $S_1$  are used to reconstruct the sentence vector.

*Output Feature Module* Question  $q$  is turned into a vector  $q = f_1(q)$  and then  $q$  is used to retrieve related entities from the memory pool.

$$\begin{cases} O_0 = Q_0 = q, E_0 = \phi \\ Q_{j-1} = h(Q_{j-2}, e_{j-1}), j = 2, 3, \dots \\ e_j = \arg \max_{e_k \notin E_{j-1}} p(e_k, Q_{j-1}); E_j = E_{j-1} \cup \{e_j\} \\ O_j = u(O_{j-1}, p(e_j, Q_{j-1}) * e_j) \end{cases} \quad (3)$$

At first,  $Q$  is initialized using  $q$ . In the  $j_{th}$  iteration,  $p(e_k, Q_{j-1})$  is the probability (or score) of  $e_k$  being selected to compose the feature vector for answering  $q$ . Note that every  $e$  is considered only once. In  $Q$ , we consider the entity selected in the previous iteration.  $Q$  is kept updated using  $e$  and  $p$ .

After several iterations, we use the final  $O_m$  as the output feature vector  $O$ . Note that if the  $O_*$  does not change much between iterations, we will omit the remaining loops. This early-stop strategy helps reduce the time cost.

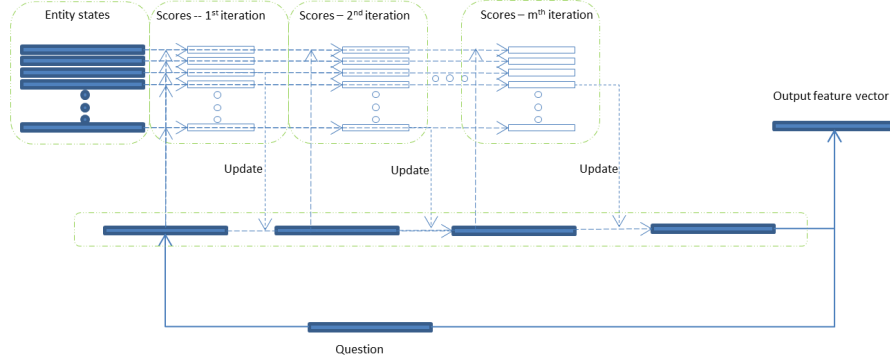


Fig. 5: the Output Feature Module: In each iteration, entities are assigned different scores which indicate their importance in constructing the output feature vector.

*Response Module* Then we decide the answer using  $a(q) = v(O)$ .  $a(q)$  produces a vector whose each item corresponds to one word in the vocabulary.  $a(q)_i$  indicates the probability of  $word_i$  being used as the correct answer. We choose the one with the highest probability. Models like recurrent neural network can be used to output a sentence as the answer.

### 3.3 Implementation

This is a supervised model and requires annotated data for the training. The training data contains the input text, questions and answers. Also we need all the entities and entities that are related to the answer labeled.

We define the function form for training as follows: As for  $f_1$ , many models, like the recurrent neural network, recursive neural network and so on [16,24,11], can be used to convert a sentence into a vector. Here we use an Long Short-Term Memory (LSTM) autoencoder [13] which takes a word sequence as input and outputs the same sequence.

$f_2$  takes a list of entity states as input and tries to reconstruct  $S_i$ . We use the Gated Recurrent Unit (GRU) [3].

$$\begin{aligned} S_i^k &= \tanh(GRU(S_i^{k-1}, e_k^i)) \\ S_i' &= S_i^j \end{aligned} \quad (4)$$

A GRU can be represented as the follows:

$$\begin{cases} z_t^j = \delta(W_z * x_t + U_z * h_{t-1})^j \\ \bar{h}_t^j = \tanh(W * x_t + U * (r_t \circ h_{t-1}))^j \\ r_t^j = \delta(W_r * x_t + U_r * h_{t-1}) \\ h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j\bar{h}_t^j \end{cases} \quad (5)$$

◦ represents an element-wise multiplication.  $z_t^j$  and  $r_t^j$  are two gates controlling the impact of historical  $h_{t-1}^j$  on the current  $h_t^j$ . The GRU takes  $x_t$  as input and updates the state of the neuron to  $h_t^j$ . Compared with LSTM which it often replaces, it simplifies the computation while still keeps a memory of previous states. Therefore it takes less time to train GRU than LSTM.

Our goal is to minimize the loss  $|S'_i - S_i|$ . Using the stochastic gradient descent, we are able to train  $f_2$  and also update  $\{e_k^i\}$ . Note that the input module and the generalization module do not interact with the remaining. Thus they can be trained in advance.

The output feature module checks the memory pool repeatedly to select entities to form a feature vector:

$$\begin{cases} Q_{j-1} = \tanh(\text{GRU}(Q_{j-2}, e_{j-1})) \\ e_j = \arg \max p(e_j, Q_{j-1}) = \arg \max \text{sigmoid}(W * \text{GRU}(e_j, Q_{j-1}) + b) \\ O_j = \tanh(\text{GRU}(O_{j-1}, p(e_j, Q_{j-1}) * e_j)) \end{cases} \quad (6)$$

To generate the final answer, we use a simple neural network which takes the feature vector  $O$  as input and predict a word as output.  $p_w = v(O) = \text{softmax}(\tanh(W' * O + b))$ . The word with the highest probability is selected. Suppose a sentence is to be generated, we use the GRU to update  $O$  and then generate the sentence  $\{w_*\}$  as follows:

$$\begin{cases} p_w^{i-1} = \text{softmax}(\tanh(W' * O_{i-1} + b)) \\ w_{i-1} = \arg \max p_w^{i-1} \\ O_i = \tanh(\text{GRU}(O_{i-1}, w_{i-1})) \end{cases} \quad (7)$$

Similar to [28], we use the stochastic gradient descent algorithm to minimize the loss function shown in Equation (8) over parameters. For an input  $S_i$  and a given question  $q$  annotated with the correct answer  $word_a$  and related entities  $\{e_r\}$ , the loss function is as follows:

$$\sum_{i \neq r} \max(0, \gamma - (p(e_r, q) - p(e_i, q))) + \sum_{l \neq a} \max(0, \gamma - (p_{word_a} - p_{word_l})) + ||\Theta||^2 \quad (8)$$

Here  $\gamma$  is the margin and  $||\Theta||^2$  is the squared sum of all parameters which is used for regularization. Note that  $\Theta$  does not include parameters of  $f_1$  and  $f_2$ . Their parameters and states of entities are learned as described in Section 2.2. Word vectors used to initialize entity states and words in autoencoder come from GloVe [17]. The dimension is set to be 50.

### 3.4 Data Preparation

The model requires entities to be annotated in advance. In this work, we treat each noun and pronoun as an entity. Different words are regarded as different entities for simplicity. This strategy saves us the effort of entity resolution which is a challenge for many languages. It also makes possible the application of the proposed model to entity



resolution <sup>4</sup>. For datasets with related entities annotated, we can use the loss function described above. But annotating the related entities is time and labour-costing. Most datasets available are not annotated. The weakly supervised learning can be applied to such data by trimming the loss function to

$$\sum_{l \neq k} \max(0, \gamma - (p_{word_k} - p_{word_l})) + ||\Theta||^2 \quad (9)$$

For unannotated data, a fully supervised training is also possible if we regard entities contained in questions as related entities or if we can use other methods to identify entities that are believed to be related.

## 4 Experiments

To verify the effectiveness of the proposed model, we conduct experiments on several datasets, including a toy QA data set bAbI [27], the large movie review dataset for sentiment classification [15] and the Machine Comprehension dataset (MC Test) [20].

### 4.1 bAbI

The example shown in Fig. 1 is extracted from the bAbI dataset. It contains 20 topics, each of which contains short stories, simple questions with regard to the stories and answers. The data is generated with a simulation which behaves like a classic text adventure game. According to some pre-set rules, stories are generated in a controlled context.

Previous work reports extremely satisfying results using memory networks for most topics (around 90% for most of them). However, we notice an interesting thing that all of them with no exception fail on the problem of path-finding which is to predict a simple path like "north, west" given the locations of several subjects. Another one is the positional reasoning. The Memory Network [27] reports accuracies of 36% and 65% for the two topics. The Dynamic Memory Network [10] reports accuracies of 35% and 60%. The proposed model (Entity-MNN) reports accuracies of 53% and 67% respectively. It is still far from satisfying but the improvements on the two tasks indicates the superiority of the entity-based memory network. For the whole dataset, we report mean error rates about 12%, comparable to 3.2 to about 24 reported by previous work [25,10,27].

The data is generated in a controlled text. As we know, QA systems trained on controlled text normally suffers when moving to real world problems [6]. Results on this toy dataset is not as convincing as that on practical tasks. Given how the bAbI data is generated, it is easy to achieve a 100% accuracy if we do simple reverse engineering to identify the entities and rules. The good results of memory networks, including our model, can not be solely attributed to their ability of comprehension. It may be partly due to their ability of inducing the entities and rules from text.

<sup>4</sup> We treat each mentions of entities as different one when processing the text and ask questions about which of these mentions refer to the same entities.

## 4.2 Machine Comprehension Test

We tested the proposed model on a dataset constructed from children stories. The machine comprehension test (MCTest) dataset [20] has 500 stories and 2000 questions (MC500). All of them are multiple choice reading comprehension questions. An additional smaller dataset with 160 stories and 640 questions (MC160) is also included in the MCTest data and used in our work.

Since the proposed model does not consider the form of multiple choice questions, we need to convert MCTest data into suitable formats firstly. When answering a multiple choice question, one is provided with several alternatives of which at least one is correct. These alternatives can be regarded as information known.

For a question, we replace the “Wh-” words using each alternative and Each alternative is turned to a new declarative sentences. These generated declarative sentences are generally understandable though may not be grammatically correct. Then we use the proposed system to decide whether the generated sentences are correct or wrong. However, we do not distinguish between questions with only one answer and those with more than one answers as these newly generated sentences are treated separately. In other words, all questions are treated as having multiple answers.

The MCTest contains only hundreds of stories and is usually used for test only as statistical models normally require a large amount of training data. However, we still obtain satisfying results using this dataset. Table 1 demonstrates the effectiveness of the entity-based model on the MCTest dataset. We outperform the previous state-of-the-art [26,22] on both MC160 and MC500. Our model does not employ rich semantic features as others do, and hence is easy to be migrated to languages aside from English.

| Sys.               | Acc.(%) MC160 |          |         | Acc.(%) MC500 |          |         |
|--------------------|---------------|----------|---------|---------------|----------|---------|
| Type               | Single        | Multiple | Average | Single        | Multiple | Average |
| Richardson’13 [20] | 76.8          | 62.5     | 69.2    | 68.0          | 59.5     | 63.3    |
| Wang’15 [26]       | 84.2          | 67.9     | 75.3    | 72.1          | 67.9     | 69.9    |
| Sachan’16 [22]     | -             | -        | -       | 72.0          | 68.9     | 70.3    |
| EntityMNN          | Average=76.1  |          |         | Average=76.6  |          |         |

Table 1: Results on Machine Comprehension Test

## 4.3 Large Movie Review Dataset

We further tested our model on the Large Movie Review Dataset [15], which is a collection of 50,000 reviews from IMDB, about 30 reviews per movie. Each review is assigned a score from 1 (very negative) to 10 (very positive). The ratio of positive samples to negative samples is 50:50. Following the previous work [15], we only consider polarized samples with scores no greater than 4 or no smaller than 7.

For each review, we present it as a short story and then add a question “what is the opinion?”. The answer is either “negative” or “positive”. In this way we turn this task into a question answering problem. Note that although here the answer to a question is either “negative” or “positive”, we do not put any constraints on the output. It is treated in the same way as open domain question answering and the system is expected to learn to predict the output by itself.

| Sys. | Maas' 11 [15] | Johnson' 14 [8] | Johnson' 15 [9] | EntityMNN % |
|------|---------------|-----------------|-----------------|-------------|
| Acc. | 89            | 93.4            | 95              | 97.2        |

Table 2: Results on Large Movie Dataset

We do not use the full dataset as the training takes a long time. We randomly select 10K samples (5K negative + 5K positive) for training and another 10K for test. We obtain an accuracy of 97.2% on the subset which is higher than previous work [15,8,9] as is shown in Table 2. By exploring relations between entities, we consider information that is usually not included for classification tasks and obtain better results.

#### 4.4 Analysis

The proposed model is designed based on the assumption that entities are the core of text. By updating the states of entities, information carried by text is encoded into entities. Thus all questions which are related to the text can be answered based on entities solely.

Using entities enable us to break a sentence into smaller text units and analyze text from a smaller scale. As stated, if one entity  $e_i$  in sentence  $S_a$  interacts with another entity  $e_j$  in sentence  $S_b$ , dealing with  $e_i$  and  $e_j$  directly is much easier than dealing with  $S_a$  and  $S_b$ . Therefore the proposed model overcomes this problem as has been proven in our experiments. A shortcoming with the proposed model is that, it cannot handle text that contains very few entities. Also hidden entities are not considered. As we know, pro-drop languages, like Japanese and Chinese, tend to omit certain classes of pronouns when they are inferable. The proposed model will encounter problems when dealing with such text.

## 5 Conclusion

This work presents the entity-based memory network model for text comprehension. All the information conveyed by text is encoded into the states of its containing entities and questions regarded to the text are answered using these entities. Experiments on several tasks have proven the effectiveness of the proposed model. The proposed model is based on the assumption that entities can express all the information of text. In future research, we will further explore its ability by considering more components in text.

## References

1. Bordes, A., Usunier, N., Chopra, S., Weston, J.: Large-scale simple question answering with memory networks. arXiv preprint arXiv:1506.02075 (2015)
2. Brill, E., Lin, J.J., Banko, M., Dumais, S.T., Ng, A.Y., et al.: Data-intensive question answering. In: TREC. vol. 56, p. 90 (2001)
3. Chung, J., Gulcehre, C., Cho, K., Bengio, Y.: Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555 (2014)

4. Dong, L., Wei, F., Zhou, M., Xu, K.: Question answering over freebase with multi-column convolutional neural networks. In: *ACL* (1). pp. 260–269 (2015)
5. Graves, A., Wayne, G., Danihelka, I.: Neural turing machines. *arXiv preprint arXiv:1410.5401* (2014)
6. Hermann, K.M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., Blunsom, P.: Teaching machines to read and comprehend. In: *Advances in Neural Information Processing Systems*. pp. 1693–1701 (2015)
7. Iyyer, M., Boyd-Graber, J.L., Claudino, L.M.B., Socher, R., Daumé III, H.: A neural network for factoid question answering over paragraphs. In: *EMNLP*. pp. 633–644 (2014)
8. Johnson, R., Zhang, T.: Effective use of word order for text categorization with convolutional neural networks. *arXiv preprint arXiv:1412.1058* (2014)
9. Johnson, R., Zhang, T.: Semi-supervised convolutional neural networks for text categorization via region embedding. In: *NIPS*. pp. 919–927 (2015)
10. Kumar, A., Irsoy, O., Su, J., Bradbury, J., English, R., Pierce, B., Ondruska, P., Gulrajani, I., Socher, R.: Ask me anything: Dynamic memory networks for natural language processing. *arXiv preprint arXiv:1506.07285* (2015)
11. Le, Q.V., Mikolov, T.: Distributed representations of sentences and documents. In: *ICML*. vol. 14, pp. 1188–1196 (2014)
12. Lehnert, W.G.: *The process of question answering: A computer simulation of cognition*. Lawrence Erlbaum Associates (1978)
13. Li, J., Luong, M.T., Jurafsky, D.: A hierarchical neural autoencoder for paragraphs and documents. *arXiv preprint arXiv:1506.01057* (2015)
14. Lin, D., Pantel, P.: Discovery of inference rules for question-answering. *Natural Language Engineering* 7(04), 343–360 (2001)
15. Maas, A.L., Daly, R.E., Pham, P.T., Huang, D., Ng, A.Y., Potts, C.: Learning word vectors for sentiment analysis. In: *ACL-HLT*. pp. 142–150 (2011)
16. Mikolov, T., Karafiát, M., Burget, L., Cernocký, J., Khudanpur, S.: Recurrent neural network based language model. In: *Interspeech*. vol. 2, p. 3 (2010)
17. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: *EMNLP*. vol. 14, pp. 1532–43 (2014)
18. Poon, H., Christensen, J., Domingos, P., Etzioni, O., Hoffmann, R., Kiddon, C., Lin, T., Ling, X., Ritter, A., Schoenmackers, S., et al.: Machine reading at the university of washington. In: *NAACL-HLT 2010 Workshop*. pp. 87–95. *ACL* (2010)
19. Ravichandran, D., Hovy, E.: Learning surface text patterns for a question answering system. In: *Proceedings of the 40th annual meeting on association for computational linguistics*. pp. 41–47. Association for Computational Linguistics (2002)
20. Richardson, M., Burges, C.J., Renshaw, E.: Mctest: A challenge dataset for the open-domain machine comprehension of text. In: *EMNLP*. vol. 3, p. 4 (2013)
21. Riloff, E., Thelen, M.: A rule-based question answering system for reading comprehension tests. In: *ANLP/NAACL2000 Workshop*. pp. 13–19. *ACL* (2000)
22. Sachan, M., Xing, E.P.: Machine comprehension using rich semantic representations. In: *ACL* (2016)
23. Shen, D., Lapata, M.: Using semantic roles to improve question answering. In: *EMNLP-CoNLL*. pp. 12–21 (2007)
24. Socher, R., Lin, C.C., Manning, C., Ng, A.Y.: Parsing natural scenes and natural language with recursive neural networks. In: *ICML*. pp. 129–136 (2011)
25. Sukhbaatar, S., Weston, J., Fergus, R., et al.: End-to-end memory networks. In: *NIPS*. pp. 2440–2448 (2015)
26. Wang, H., McAllester, M.B.K.G.D.: Machine comprehension with syntax, frames, and semantics. *ACL, Volume 2: Short Papers* p. 700 (2015)

27. Weston, J., Bordes, A., Chopra, S., Rush, A.M., van Merriënboer, B., Joulin, A., Mikolov, T.: Towards ai-complete question answering: A set of prerequisite toy tasks. arXiv preprint arXiv:1502.05698 (2015)
28. Weston, J., Chopra, S., Bordes, A.: Memory networks. arXiv preprint arXiv:1410.3916 (2014)